

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-318766

(43)Date of publication of application : 16.11.2001

(51)Int.Cl.

G06F 3/06

G06F 12/08

G06F 12/16

G11B 19/02

(21)Application number : 2000-135013

(71)Applicant : NEC SOFTWARE SHIKOKU LTD

(22)Date of filing : 08.05.2000

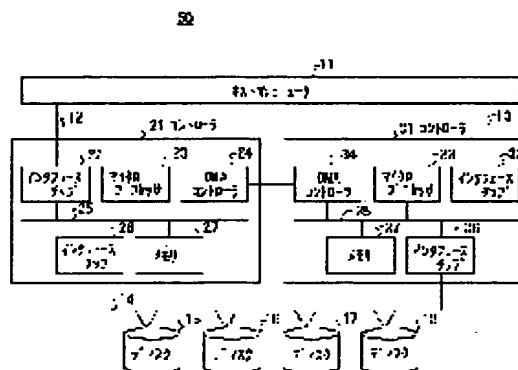
(72)Inventor : SHIRAISHI KAZUYA

(54) DISK ARRAY DEVICE AND CACHE MEMORY CONTROL METHOD

(57)Abstract:

PROBLEM TO BE SOLVED: To provide a disk array device and a cache memory control method for executing second and succeeding pieces of write command processing without putting loads on the processor of a main controller.

SOLUTION: This disk array device is composed of two controllers and executes the double write of write data in the two controllers without using the hardware of a shared memory. The write data are held in both of the two controllers and a means is provided for making the write data held in the other one of the two controllers mutually referable when a fault is generated in one of the controllers.



12、15、16、17は、
各コントローラ内の
キャッシュメモリ、
プロセッサ、および
ディスクアレイ。

LEGAL STATUS

[Date of request for examination] 17.04.2001

[Date of sending the examiner's decision of rejection] 12.02.2003

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

Best Available Copy

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

た基板（パッケージ）で作成されているため、当該メモリの故障を考慮すると、共有メモリ3000の2重化が必要になり、回路構成が大きくコストが高くなるという欠点があり、中小型システムでは採用することが難しいという問題点があった。

【0018】また、上記第2従来技術は、内部バスをオリジナルバスで構成することになるため、汎用バスを採用できる場合はインタフェースチップなどを汎用品の中から自由に選択でき、高性能な装置を安価に入手できるメリットがあるものの、オリジナルバスではこのような構成が難しいという問題点があった。

【0019】そして、上記第3従来技術は、2重書きにはディスク側の汎用バス14を使用するため、処理のオーバーヘッドが大さいという問題点があり、また、メインコントローラのプロセッサに負荷がかかるという問題点もあった。

【0020】本発明は斯かる同趣点を覆みてなされたものであり、その目的とするところは、メイトコントローラのセグメントをシーズしてそのキャッシュアドレスとビットマップを保持することにより、2個目以降のライトコマンド処理においてメイトコントローラのプロセッサに負荷をかけることなく実行できるようにするディスクレイアウトおよびキャッシュメモリ制御方法を提供する点にある。

[0021]

【課題を解決するための手段】この発明の請求項1に記載の発明の要旨は、2つのコントローラで構成されかつ共有メモリのハードウェアを用いることなく該2つのコントローラ間でライトデータの2番書きを実行するディスクレイアウト装置において、ライトデータと前記2つのコントローラの両方に保持するとともに、前記コントローラの一方に故障が発生した後に前記2つのコントローラの他方に保持されているライトデータを相互に複製可能とする手段を有することを特徴とするディスクレイアウト装置にある。また、この発明の請求項2に記載の発明の要旨は、ホストコンピュータとディスクとの間に接続2つのコントローラを備え、前記ホストコンピュータと前記2つのコントローラの両方が共用バスで接続されるとともに、前記2つのコントローラの両方と前記ディスクが両方に接続されるとともに特徴とする請求項1に記載のディスクレイアウト装置に存する。また、この発明の請求項3に記載の発明の要旨は、前記2つのコントローラのそれぞれは、マクロプロセッサと、リードまたはライトのデータをキャッシュしかつプロセッサの制御動作を制御するメモリ、前記ホストコンピュータとの通信を制御する第1のインターフェースチップと、前記ディスクとの通信を制御する第2のインターフェースチップと、前記メモリに任意でアドレスのデータを前記2つのコントローラの他方の前記メモリに任意アドレスに逐逐で転送DMAコントローラを有していることを特徴

する請求項４および５に記載のディスクアドレスに対応する位置にある。またこの発明の請求項４に記載の発明の発行者は、セグメント単位でディスクアドレスを管理する前記ステップには、セグメント管理情報と、ビットマップと、ライトキャッシュが定義されることを特徴とする請求項３に記載のディスクアドレスに該当する。また、この発明の請求項５に記載の発明の発行者は、前記セグメント管理情報はセグメントの状態を保持するものであるため、シーリング処理を示すディレクトリと、ロジカルコンテントリストと、ロジカルブロックアドレスと、コントローラの両方の前に記述されたビットマップを有する自己ビットマップアドレスと、コントローラへのキャッシュアクセスを示す自己ビットマップアドレスと、前記２つのコントローラ間のデータの送受信を行う前記ビットマップ格納位置を示す他ビットマップアドレスと、前記２つのコントローラの両方のキャッシュ位置を示す他キャッシュアドレスを有することを特徴とする請求項４に記載のディスクアドレスに該当する。また、この発明の請求項６に記載の発明の発行者は、前記ビットマップは、セグメントに対応する各キャッシュのどの位置に有効なデータが存在するかを示すビットマップであることを特徴とする請求項５に記載のディスクアドレスに該当する。また、この発明の請求項７に記載の発明の発行者は、前記ライトキャッシュは、各セグメントに対応するキャッシュスキューモータによってアクセスされることを特徴とする請求項６に記載のディスクアドレスに該当する。また、この発明の請求項８に記載の発明の発行者は、前記ホストコンピュータからのライトコマンドを受けた後に前記コントローラの方は自己の前記セグメント管理情報を参照して自身がサイズとして指定したバッファサイズを超えて前記２つのコントローラ間のデータを送信または受領してセグメントをシージングし、前記ホストコンピュータからライトデータを受領して自己の前記ライトキャッシュに格納した後該格納位置に対応した自己の前記ビットマップを有効状態に更新し、自己の前記ライトキャッシュのデータおよび自己の前記ビットマップを前記２つのコントローラ間の他の前記DMAコントローラを介して前記２つのコントローラ間の他の前記ライトキャッシュおよび前記ビットマップに送信してライトコマンドを受け渡し、前記コントローラの方でシージングを行いセグメントにヒットした場合に前記２つのコントローラの間方にシージングコマンドを送信することなく、前記ホストコンピュータからライトデータを受領して自己の前記ライトキャッシュに格納した後該格納位置に対応した自己の前記ビットマップを有効状態に更新し、自己の前記ライトキャッシュのデータおよび自己の前記ビットマップを前記２つのコントローラ間の他の前記DMAコントローラを介して前記２つのコントローラ間の他の前記ライトキャッシュと前記ビットマップに送信することを特徴とする請求項７に記載のディスクアドレスに該当する。また、この発明の請求項９に係る請求項

記の発明の要旨は、２つのコントローラで構成されたシステムにおいてハードウェアを用いることなく各該２つのコントローラでライブラリデータの転送並行を実行するイスクエリと答復に対して、ライブラリデータを前記２つのコントローラの両方に保持するとともに、前記コントローラの方から読取されるデータに前記２つのコントローラの方から提供されているライブラリデータを相対し得る可とする工程を有することと特徴とするキャッシュメモリ制御方法に在る。また、この発明の請求項１０に記載の発明の要旨は、前記２つのコントローラそれぞれから受け取られるデータまたはライブラリのデータをセグメント単位でキャッシュしかつフロッパの制御機能を兼充するメモリに、セグメント管理情報と、ビットマップと、自己キャッシュを定義する工程を有することと特徴とし、その請求項９に記載のキャッシュメモリ制御方法に在る。また、この発明の請求項１１に記載の発明の要旨は、前記セグメント管理情報は各セグメントの状態を示すものであるとして、サイズの値を表示するサイズワードと、ロジカルユニット番号と、ロジカルブロックアドレスと、自コントローラの両方のビットマップ格納位置を示す自己ビットマップアドレスと、自コントローラのキャッシュ位置を示す自己キャッシュアドレスと、前記２つのコントローラの方のビットマップ格納位置を示す他のビットマップアドレスと、前記２つのコントローラの方のキャッシュ位置を示す他キャッシュアドレスを有することと特徴とする請求項１０に記載のキャッシュメモリ制御方法に在る。また、この発明の請求項１２に記載の発明の要旨は、前記ビットマップは、各セグメントに対応するキャッシュとの位置に有効データが存在することを示すビットマップであることを特徴とする請求項１１に記載のキャッシュメモリ制御方法に在る。またこの発明の請求項１３に記載の発明の要旨は、前記ビットマップは、色セグメントに対応する請求項１２に記載のキャッシュメモリ制御方法に在る。また、この発明の請求項１４に記載の発明の要旨は、前記バスコネクティブからイトコンポートを受信した後に前記コントローラの方には、自分の前記セグメント管理情報を参照して自身が一貫しているセグメントにヒットしているかを判断する工程と、ヒットしていなければ別ポートを介して前記２つのコントローラの方からサイズワードを送信してセグメントをサイズする工程と、前記バスコネクティブから入力されたデータを受けて自分の前記ライブラリキャッシュに格納し該格納位置に対応した自分の前記ビットマップの有効データを更新する工程と、自分の前記ビットマップのデータに基づきより自分の前記ビットマップを前記２つのコントローラの方のDMAコントロールを介して前記２つのコントローラの方の格納ライブラリキャッシュおよび前記ビットマップに送信してライブラリコマンドを受信する工程と、前記２つのコントローラの方から自己

ていてメモリバンクにヒットした場合には前記２つのコントローラのうち一方のメモリバンクにコマンドを送信することなく、前記コントローラにコンテスタから格納されたデータを読み出し、自己の格納したメモリーバンクに格納し、その後格納装置に自己の格納した自己の前記ビットマップを有効状態に更新する工程と、自己の前記メモリーバンクのデータおよび自己の前記ビットマップを前記２つのコントローラの一方のDMAコントローラを介して前記２つのコントローラの一方の格納したメモリーバンクと前記ビットマップとに送付する工程を有することと特徴とする請求項１に記載のキャッシュメモリ制御方法に存する。

【００２２】

【発明の実施の形態】 本発明は、２コントローラ構成で、かつ、各メモリーバンクのハードウェアを１つのメモリーバンクに分割するものにおいて、ライトデータを格納する際の２重書きに際して、メモコントローラ（後述する本発明の相手コントローラ）のプロセッサに負荷をかけることなくと特徴としている。

【００２３】 ライトデータはコントローラ内のメモリに格納されるため、１コントローラ構成で自己のメモリーバンクを使用すると、コントローラの故障によってライトデータが失ってしまう。一方、２コントローラ構成で、かつ、ライトデータを両方のコントローラに持つておけば、片方のコントローラが故障しても、ライトデータが失われることはない。このようにライトデータを両方のコントローラに保持する仕組みを二重書きと呼ぶ。以下、本発明の実施の形態を図面に基づいて詳細に説明す

【0024】（第1の実施の形態）図1は、本発明の第1の実施の形態に係るデータスクリーン装置50を説明するためのブロック図である。本実施の形態のデータスクリーン装置50は、ホストコンピュータ11とデータスクリーン15、…、18との間に2つのコントローラ21、31を有するものである。

【0025】ホストコンピュータ11とコントローラ21、31は、例えば、SCS1（スカジー・スモール・コンピュータ・システム・インターフェース）バスのような汎用バス22、13で接続されている。

【0026】コントローラ21、31はデータスクリーン15、…、18は、例えば、SCS1バスのような汎用バス14で接続されている。

【0027】コントローラ21は、マイクロプロセッサ23と、リードまたはライトのデータをキャッチし、かつプロセッサの制御部を記憶するメモリ24と、ホストコンピュータ11との通信を制御するインターフェース部25と、データバス15、…、18との通信を制御するインターフェース部26と、自身のコントローラ21のメモリ27の任意アドレスのデータをホストコンピュータのメモリ27の任意アドレスに送受するDMAコントローラ24を有し、コントローラ21内の各

トップ内部バス25で接続されている。

[0028] 図1は、コントローラ31は、マイクロプロセッサ33と、リードまたはライトのデータをキャプチャングおよびプロセスする制御機能を含むメモリ37と、バストコンピュータ11との通信を制御するインタフェースチップ32と、ディスク15、…、18との通信を制御するインタフェースチップ34、…、18とのコントローラ31のメモリ37の任意アドレスのデータをメインコントローラのメモリ37の任意アドレスに送受できるDMAコントローラ34と、コントローラ33内のキャッシュは内部バス35で接続されている。

[0029] 図2は、図1のメモリ37、37に記憶される情報の最上位である。本明記ではセグメントと呼ばれる単位でキャッシュを管理しており、メモリ27、37には、セグメント管理情報41、44と、ビットマップ42、45と、ライトキャッシュ43、46が記憶されている。セグメント管理情報41、44は、各セグメントの属性を保持するもので、サイズの属性を各セグメントフラグと、LUN（ロジカルユニット番号）と、LBA（ロジカルブロックアドレス）と、自コントローラのビットマップ格納位置を示すビットマップアドレスと、自コントローラのキャッシュ位置を示す自キャッシュアドレスと、メインコントローラのビットマップ格納位置を示す他ビットマップアドレスと、メインコントローラのキャッシュ位置を示す他キャッシュアドレスを有する。ビットマップ42、45は、各セグメントに対応するキャッシュのその位置に有効なデータが存在するかを示すマップで、0が無効、1が有効である。ライトキャッシュ43、46は、各セグメントに対応するキャッシュメモリ自体である。

[0030] コントローラ2がバストコンピュータ1からライトデータを受ける時、セグメント管理情報41を参照して、自身がサイズしているセグメントにヒットしているかを確認し、ヒットしていなければ読み出し46をしてコントローラ33にキャッシュデータを送渡し、セグメントをシーズする。その後、ライトコンピュータ1からライトデータを受けたライトキャッシュ44に格納し、当該格納位置に格納したビットマップ42と1に更新する。その後、ライトキャッシュ43のデータとビットマップ44をDMAコントローラ2、43を介してコントローラ31のメモリ37のライトキャッシュ46にビットマップ45に送渡する。この処理で、セグメントをシーズするの処理によって、コントローラ31のライトキャッシュ33が格納されて

チップ42と接続する。その後、ライトキャッシュ3のデータとビットマップ42とDMAコントローラ24、34を介してコントローラ31のライトキャッシュ46とビットマップ45と送受信する。

【0030】以上説明したように第1の実施の形態によれば、2回以降のライトコマンドでコントローラ31のプロセッサ33に負荷がかかることなくライトコマンド処理が実行できるため、装置全体の性能が向上する。

【0031】次にディスクレイアウト装置60の動作(キャッシュセクタ制御方法)について説明する。なお、本実施の形態の動作について、ホストコンピュータ11からライトコマンドを受信し、それを2番書きする例を示す。また、本実施の形態では2回以降のライトコマンド処理に負荷が表れるため、ライトコマンドを連続して2回受信した場合で説明する。

【0034】図3は、本発明の第1の実施の形態に係るキャッシュセクタ制御方法におけるライト処理を説明するためのフローチャートである。実行パス14にSCSインターフェース4を用い、ライトコマンドを2回受信した場合の処理の流れを示している。

【0035】コントローラ2は、ホストコンピュータ11からライトコマンドを受信し(ステップ55:ライトコマンド受信)、自身のセグメント管理情報44に参照してセグメントのビットマップの判定を行う(ステップ56:ビットマップ判定)。1回目のコマンドなのでそれがミスし、コントローラ31に汎用バス14を介してセグメント管理44を出力し(ステップ65:SCS1でセグメント管理44出力)、同時にそのバリエータとして、UNとLBAを送信する。

【0036】それを受信したコントローラ31は、自身の空きセグメントをセイズして、セグメント管理情報44に、セイズ状態を示すフラグとUNとLBAとビットマップ42に、セイズしたアドレスとそれに対応するビットマップ42、45を格納するアドレスを登録して(ステップ54:セグメント管理情報44に登録)、コントローラ2にシイズの完了報告を出力し(ステップ65:SCS1でシイズ完了報告)、同時にそのレスポンス情報としてキャッシュアドレスとビットマップ格納アドレスを送信する。コントローラ2は自身のセグメント管理情報44に、セイズ状態を示すフラグとUNとLBAとLBA自身がセイズしたキャッシュのアドレスとそれに対応するビットマップ42、45を格納するアドレスとそれに対応するビットマップ42、45を格納するアドレスを登録する(ステップ54:セグメント管理情報44に登録)。

【0037】その後、ホストコンピュータ11から自身がセイズしているキャッシュアドレスにライトデータを出力し(ステップ65:7:ライトデータ出力)、それに

に対応するビットマップ42を更新する(ステップS56)。ビットマップ更新)。

[0038] 同実行データで、DMAコントロール44、34を先入れコントロール31からシークしているキャッシュアドレスにダイレクトに転送し(ステップS59: DMAでビットデータ転送)、同ビットマップ42、45を、同じくDMAコントロール44、34を先入れコントロール31のビットマップ45に転送し(ステップS60: DMAでビットマップ転送)、これによって2重書きが完了し同時に1回目のライトコマンドが完了する。

[0039] 続けて、2回目のライトコマンドを受理する(ステップS61: ライトコマンド受付)。自身セグメント管理情報44を参照してセグメントのビットマップの特定を行い(ステップS62: ビットマップ特定)、それが終了すると、1回目のより大きなサイズの取り無しにホストコンピュータ11から自身がシークしているキャッシュアドレスにライトデータを受信することができる(ステップS63: ライトデータ受付)。

[0040] その後、それらに対応するビットマップ44を更新し(ステップS64: ビットマップ更新)、ライトデータと同ビットマップ44、45を先入れ1回目と同様にコントロール31に転送し(ステップS65: DMAでライトデータ転送)、ステップS66: DMAでビットマップ転送)、これによって2重書きが完了し同時に2回目のライトコマンドが完了する。

[0041] これらの処理を比較すると、1回目のコマンド処理(ステップS67)では、サイズコマンドの取り取りが行われる処理に長時間がかかっているのに対し、2回目以降のコマンド処理(ステップS68)ではサイズコマンドの取り取りがなく、処理時間が短縮されるので、コントロール31のプロセッサに負荷がかかっていないことがわかる。

[0042] 図4は、本発明の第1の実施形態に係るディスクライブラリ装置50の各コントロール21、31におけるライト処理を説明するためのフローチャートである。コントロール21とコントロール31のライトコマンド処理のフローチャートを示している。

[0043] 図4を参照してそれぞれのプロセスがどのような処理を行うかを説明する。ステップ71(コントロール21の処理)の実行後、コントロール21は、ホストコンピュータ11からライトコマンドを受信すると(ステップS72: ホストコンピュータ11からライトコマンド受付)、自身からシークしているセグメントにヒットしているかどうかを判断する(ステップS73: シークとセグメントにヒット)。

[0044] ここでヒットならミスするわけで(ステップS73のNo)、SCS1にヒットしたコントロール31セグメントでミスデータを転送する(ステップS74: SCS1にコントロール31のミスデータ転送)。

(下)行)。

[00045] コントローラ31はそのシードコマンドを受信(ステップS82: コントローラ31の処理/セグメント管理情報4を受信)。SCS1にてコントローラ2からセグメントシードコマンドを受信し、自身のセグメント管理情報44を更新して(ステップS84: セグメント管理情報44の更新)、コントローラ2にてセグメントシード成功の通知を送信する(ステップS85: SCS1にてコントローラ2へセグメントシード成功の通知を送信)。

[00046] コントローラ21はその通知を受信し(ステップS75: SCS1にてコントローラ31からセグメントシード成功の通知を受信)、自身のセグメント管理情報41を更新して(ステップS76: セグメント管理情報41の更新)、ホストコンピュータ11からライトデータを受信し(ステップS77: ライトコンピュータ11からライトデータを受信)、それに対応してビットマップ42を更新する(ステップS78: ビットマップ42の更新)。

[00047] その後、DMAにてコントローラ31にライトデータを送信し(ステップS79: DMAにてコントローラ31にライトデータを送信)、続いてビットマップ42を送信し(ステップS80: DMAにてコントローラ31へビットマップ42を送信)、これらの実行をもってコントローラ31へ完了報告を行う(ステップS81: ホストコンピュータ11へ完了報告)。

[00048] 2回目以降は、シード中セグメントのヒット判定(ステップS73: シード中セグメントにヒット)に代ってヒットとなく(ステップS73のYes/No)、シード関連の処理(ステップS74: SCS1にてコントローラ31へセグメントシードコマンドを送信、ステップS76: SCS1にてコントローラ31からセグメントシード成功の通知を受信、ステップS76: セグメント管理情報41の更新)をスキップする。これにより、コントローラ31の処理/セグメント管理47が行われることなく、マイクロプロセッサ33の負荷がかけられない。

[00049] 以上説明したように第1の実施の形態によれば、メインコントローラのセグメントをローゼンしてそのキャッシュアドレスにビットマップ41、45を保持することにより、2回以上繰り返すライトコマンド処理においてはメインコントローラのプロセッサに負荷をおけることなく実行できることになり、その結果、装置全体の性能の向上に寄与することができるようになるという効果を得る。

[00050] (第2の実施の形態) 以下、本発明の第2の実施の形態を説明する。なお、本発明1の実施の形態において既に述べたものと同様の部分については、同一符号を用い、重複した説明は省略する。

[00051] 本発明の第2の実施の形態は、上記第1の

実装の形態の基本構成に加えて、コントローラ21、31を2増やすこともできる構成を備えている点に特徴を有している。具体的には、シーズコマンドを送信するための、例えば、SCSIバスのような汎用バスと、自身のメモリ27、37のデータを送信するためのコントローラ21、31のメモリ27、37にダイレクトに接続できるDMAコントローラ（本図示）を有している。

【0052】処理の流れは、上記第1の実装の形態の2コントローラ構成の場合と基本的に同じであり、2番目のシーズコマンドを送信してセグメントをシーズし、2回目以降はシーズ処理無しで2重書きを実行する。

【0053】なお、本発明が上記各実施の形態に限定されず、本発明の技術思想の範囲内において、各実施の形態は適宜変更され得ることは明らかである。また上記構成要素の数、位置、形状等は上記実施の形態に限定されず、本発明を実施する上で好適な数、位置、形状等にする事ができる。また、各図において、同一構成要素には同一符号を付している。

【0054】
【発明の効果】本発明は以上のように構成されているので、メイトコントローラのセグメントをシーズしてそのキャッシュアドレスとビットマップを保持することにより、2回目以降のライトコマンド処理においてはメイトコントローラのプロセッサに負荷をかけることなく実行できるようになり、その結果、装置全体の性能の向上を図ることができるものとなるという効果を奏する。
【図面の簡単な説明】
【図1】本発明の第1の実装の形態に係るディスクレイ

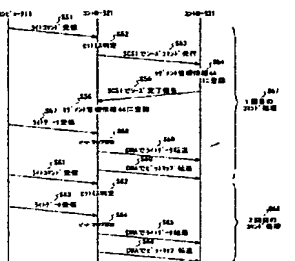
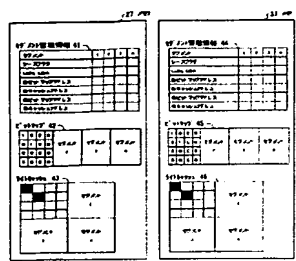
装置を説明するための機能ブロック図である。
【図2】図1のメモリに記憶する情報の定義図である。
【図3】本発明の第1の実装の形態に係るキャッシュメモリ制御方法におけるライト処理を説明するためのフローチャートである。
【図4】本発明の第1の実装の形態に係るディスクレイ装置の各コントローラにおけるライト処理を説明するためのフローチャートである。
【図5】第1従来技術のキャッシュメモリ制御方法を説明するための機能ブロック図である。
【図6】第2従来技術の主記憶装置の制御方法を説明するための機能ブロック図である。
【符号の説明】
11…ホストコンピュータ
12、13、14…汎用バス
15、…、18…ディスク
21、31…コントローラ
22、26、32、36…インタフェースチップ
23、33…マイクロプロセッサ
24、34…DMAコントローラ
25、35…内部バス
27、37…メモリ
41、44…セグメント管理情報
42、45…ビットマップ
43、46…ライトキャッシュ
50…ディスクレイ装置

【図1】本発明の第1の実装の形態に係るディスクレイ装置を説明するための機能ブロック図である。
【図2】図1のメモリに記憶する情報の定義図である。
【図3】本発明の第1の実装の形態に係るキャッシュメモリ制御方法におけるライト処理を説明するためのフローチャートである。
【図4】本発明の第1の実装の形態に係るディスクレイ装置の各コントローラにおけるライト処理を説明するためのフローチャートである。
【図5】第1従来技術のキャッシュメモリ制御方法を説明するための機能ブロック図である。
【図6】第2従来技術の主記憶装置の制御方法を説明するための機能ブロック図である。
【符号の説明】
11…ホストコンピュータ
12、13、14…汎用バス
15、…、18…ディスク
21、31…コントローラ
22、26、32、36…インタフェースチップ
23、33…マイクロプロセッサ
24、34…DMAコントローラ
25、35…内部バス
27、37…メモリ
41、44…セグメント管理情報
42、45…ビットマップ
43、46…ライトキャッシュ
50…ディスクレイ装置

【図1】本発明の第1の実装の形態に係るディスクレイ装置を説明するための機能ブロック図である。
【図2】図1のメモリに記憶する情報の定義図である。
【図3】本発明の第1の実装の形態に係るキャッシュメモリ制御方法におけるライト処理を説明するためのフローチャートである。
【図4】本発明の第1の実装の形態に係るディスクレイ装置の各コントローラにおけるライト処理を説明するためのフローチャートである。
【図5】第1従来技術のキャッシュメモリ制御方法を説明するための機能ブロック図である。
【図6】第2従来技術の主記憶装置の制御方法を説明するための機能ブロック図である。
【符号の説明】
11…ホストコンピュータ
12、13、14…汎用バス
15、…、18…ディスク
21、31…コントローラ
22、26、32、36…インタフェースチップ
23、33…マイクロプロセッサ
24、34…DMAコントローラ
25、35…内部バス
27、37…メモリ
41、44…セグメント管理情報
42、45…ビットマップ
43、46…ライトキャッシュ
50…ディスクレイ装置

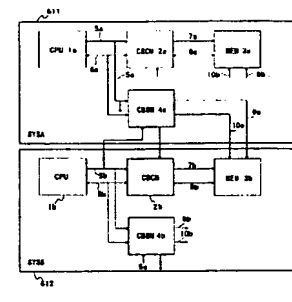
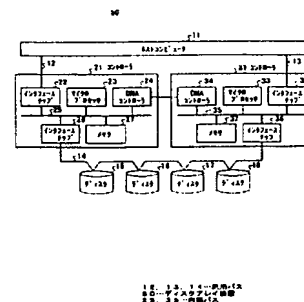
【図2】

【図3】

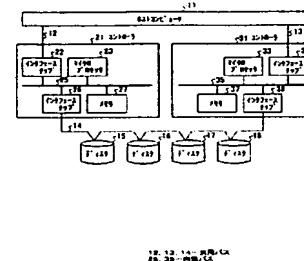


【図1】

【図6】



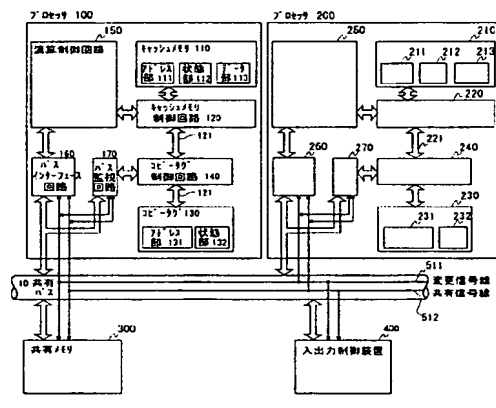
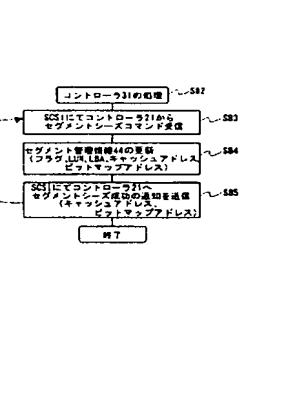
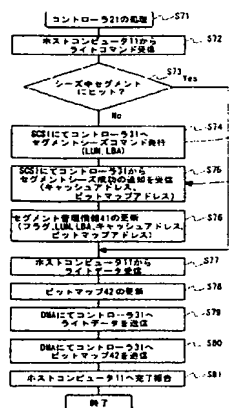
【図7】



12、13、14…汎用バス
15、16、17、18…ディスク

【図4】

【図5】



フロントページの続き

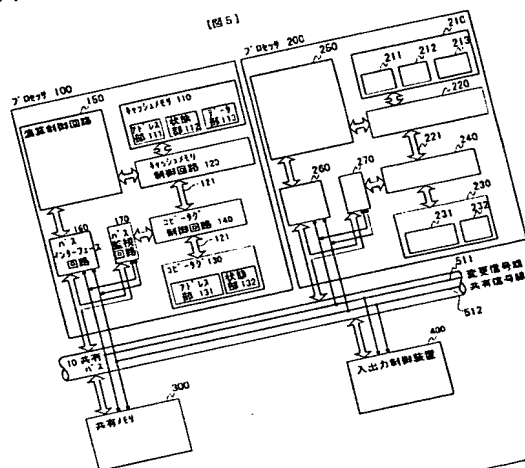
(St) 1st Cl.	発明記号	F I	備考
G 0 6 F 12/16	3 1 0	G 0 6 F 12/16	3 1 0 J
	3 2 0		3 2 0 I
G 1 1 B 19/02	5 0 1	G 1 1 B 19/02	5 0 1 D



2003 04 24 10.41

- 10 -

ディスクアレイ装置およびキャッシュメモリ制
御方法



(51) Int. Cl.⁷
G 0 6 F 12/16
G 1 1 B 19/02

310
320
50

FI
G06F 12/16
G11B 19/02

G I I B 19/02

7-23-4 (22)

310J
320L
501F

2003 04 24 10:41

- ▶ Data that is accessed normally with some locality of reference will use partial track mode staging. This is the default mode.
- ▶ Data that is not a regular format, or where the history of access indicates that a full stage is required, will set the full track mode.
- ▶ The adaptive caching mode data is stored on disk and is reloaded at IML

Sequential reads

Cache space is released according to Least Recently Used (LRU) algorithms. Space in the cache used for sequential data is freed up quicker than other cache or record data. The ESS will continue to pre-stage sequential tracks when the last few tracks in a sequential staging group are accessed.

Stage requests for sequential operations can be performed in parallel on the RAID array, giving the ESS its high sequential throughput characteristic. Parallel operations can take place because the logical data tracks are striped across the physical data disks in the RAID array.

3.28 NVS and write operations

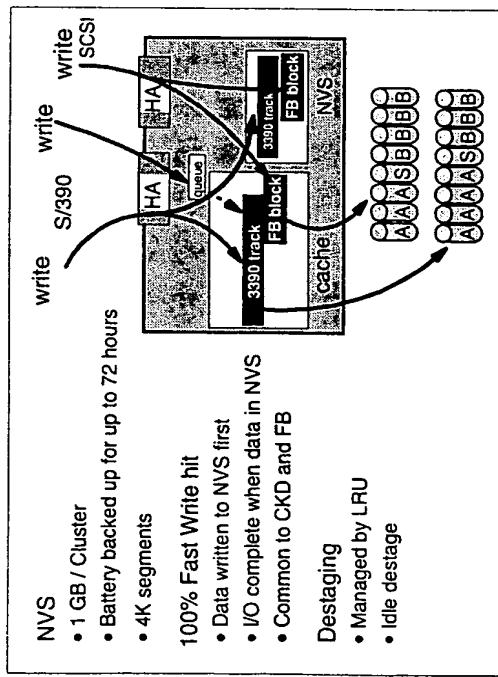


Figure 3-32 NVS - write

As Figure 3-32 illustrates, at any moment there are always two secured copies of any update into the ESS.

3.28.1 Write operations

Data written to an ESS is almost 100% fast write hits. A fast write hit occurs when the write I/O operation completes as soon as the data is in the ESS cache and non-volatile storage (NVS). The benefit of this is very fast write operations.

Fast write

Data received by the host adapter is transferred first to the NVS and a copy held in the host adapter buffer. The host is notified that the I/O operation is complete as soon as the data is in the NVS. The host adapter, once the NVS transfer is complete, then transfers the data to the cache.

The data remains in the cache and NVS until it is destaged. Destage is triggered by cache and NVS usage thresholds.

3.28.2 NVS

The NVS size is 2 GB (1 GB per cluster). The NVS is protected by a battery. The battery will power the NVS for up to 72 hours following a total power failure.

NVS LRU

NVS is managed by a Least Recently Used (LRU) algorithm. The ESS attempts to keep free space in the NVS by anticipatory destaging of tracks when the space used in NVS exceeds a threshold. In addition, if the ESS is idle for a period of time, an idle destage function will destage tracks until, after about 5 minutes, all tracks will be destaged.

Both cache and NVS operate on LRU lists. Typically space in the cache occupied by sequential data is released earlier than space occupied by data that is likely to be re-referenced. Sequential data in the NVS is destaged ahead of random data.

When destaging tracks, the ESS attempts to destage all the tracks that would make up a RAID stripe, minimizing the RAID-related activities in the SSA adapter.

NVS location

NVS for cluster 1 is located physically in I/O drawer of cluster 2, and vice versa. This ensures that we always have one good copy of data, should we have a failure in one cluster.

See 3.8, "Cluster operation: failover/failback" on page 60 for more information.

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☒ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☐ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☐ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.